

# Benchmarking Mixture-of-Recursion (MoR) on Small Language Models

**Abstract**—This paper benchmarks the Mixture-of-Recursion (MoR) architecture, which dynamically adapts computational depth for efficient language modeling. Unlike standard Transformers that apply fixed depth to all tokens, MoR uses lightweight routers to skip, execute, or recurse through layers based on input complexity. We evaluate MoR against fixed-depth Transformers on character-level next-token prediction using Tiny Shakespeare. Two experiments are conducted: (1) Efficiency profiling, where a 12-layer MoR model achieves 33% depth reduction ( $E=8$ ) while maintaining accuracy comparable to a 12-layer baseline; (2) Equal-cost comparison, where MoR ( $N=12$ ,  $E\approx 6$ ) outperforms a 6-layer baseline by 10% absolute accuracy on held-out data at equivalent computational cost. Results demonstrate that MoR’s adaptive routing enables superior efficiency and generalization, positioning it as a promising approach for resource-constrained language modeling.

**Index Terms**—Mixture-of-Recursion (MoR), LLM, SLM, Transformer, Recursive wise KV cache, Routing Mechanism (Adapted from [1]).

## I. INTRODUCTION

### A. Background of The Research

Large Language Models (LLMs) built on the Transformer architecture have established themselves as the cornerstone of modern artificial intelligence, demonstrating exceptional performance in tasks such as text generation, code synthesis, and question answering [2]–[4]. The efficacy of Transformers stems from their capacity to process sequences in parallel and capture complex patterns within extensive datasets.

However, the standard Transformer architecture suffers from significant limitations. Primarily, it imposes a fixed computational depth across all inputs, disregarding the inherent variability in token complexity. Consequently, simple tokens receive the same computational resources as complex reasoning phrases, leading to inefficient resource allocation. Furthermore, the quadratic scaling of the self-attention mechanism with respect to sequence length results in substantial memory and computational overhead [5]–[7]. These factors collectively render large models expensive to train and challenging to deploy.

This has motivated growing interest in Small Language Models (SLMs), which aim to achieve competitive performance using fewer parameters and lower computational budgets. SLMs are particularly important for deployment in resource-constrained environments such as edge devices, mobile platforms, and real-time systems [8]. However, because SLMs have limited capacity, they must utilize computation more efficiently. This makes adaptive computation mechanisms especially critical for SLM-based architectures.

Recent studies indicate that next-token prediction facilitates the learning of critical linguistic aspects, including grammar, semantics, and reasoning patterns [9], [10]. Yet, the complexity of predicting the subsequent token varies significantly. While some tokens are predictable via local context or frequency patterns, others necessitate long-range context, multi-step reasoning, or recursive processing.

To mitigate these issues, dynamic computation paradigms have emerged, enabling models to adjust layer usage based on input complexity. Mixture-of-Experts (MoE) exemplifies this approach by routing tokens to specific expert layers to conserve computation [11]. However, MoE is limited to spatial routing and lacks the capability for temporal recursion, meaning it cannot re-apply a layer to deepen reasoning.

The Mixture-of-Recursions (MoR) architecture addresses this limitation by offering each layer three execution choices: skip, execute once, or recurse. This flexibility allows for adaptive recursion, where complex tokens receive increased computation [1]. Such adaptability is vital for tasks with variable difficulty, such as code generation and question answering, positioning MoR as a potent architecture for enhancing both efficiency and reasoning capabilities.

### B. Problem Definition

Although Transformers are powerful, their fixed-depth structure creates several challenges:

- **Computational inefficiency:** Every input passes through the full layer stack, even when not necessary [12].
- **Lack of recursive computation:** Existing dynamic-routing methods like MoE cannot re-visit a layer, meaning the model cannot naturally perform deeper iterative reasoning [11].
- **Scaling bottlenecks:** Larger models demand huge memory and compute budgets, making them impractical for many applications. Continuously increasing model size to gain performance is economically and environmentally unsustainable. There is a strong need for smarter and more efficient architecture that uses less parameters.
- **Limited benchmarking of MoR:** While MoR has shown benefits on general language modeling, its performance on tasks that require deeper processing remains underexplored [1].

These limitations raise important questions:

- Can MoR reduce computational cost while maintaining accuracy?
- Does MoR perform better than a fixed-depth Transformer when both use the same amount of computation?

- How does recursion frequency correlate with input difficulty in next token prediction tasks?

To address these inquiries, this study conducts a comparative analysis of MoR and standard Transformers.

### C. Research Aim and Objectives

The primary objective of this research is to evaluate the efficacy of the Mixture of Recursions (MoR) architecture as an adaptive computation method for next-token prediction. Furthermore, we aim to benchmark MoR against standard fixed-depth Transformers to ascertain which approach offers superior efficiency and performance.

To achieve this aim, the research focuses on the following objectives:

- **Benchmark and Evaluate Performance:** Measure how well MoR and standard Transformers perform looking at accuracy and effective depth—using small-scale models.
- **Test MoR’s Efficiency:** Examine whether MoR can use less computation (lower effective depth) while still keeping the same or better performance.
- **Comparison:** Compare MoR and fixed-depth Transformers under the same computational budget for next token prediction tasks.
- **Understand MoR’s Behavior:** Study MoR’s internal routing, such as recursion rates, skip rates, and execution patterns, to see how the model adjusts its computation based on task difficulty.

### D. Contributions

This work makes the following contributions:

- **Empirical Benchmark:** We provide the first systematic evaluation of the MoR architecture on character-level language modeling using Tiny Shakespeare, demonstrating its applicability to small-scale settings.
- **Efficiency Analysis:** We show that MoR achieves a 33% reduction in effective computational depth while maintaining accuracy equivalent to fixed-depth baselines.
- **Generalization Evidence:** We demonstrate that under equivalent compute budgets, MoR achieves 10% higher held-out accuracy than shallow baselines, validating adaptive depth allocation.
- **Routing Behavior:** We analyze MoR’s routing decisions, providing insights into how dynamic recursion adapts to token-level complexity in next-token prediction tasks.

## II. LITERATURE REVIEW

### A. Introduction

Large Language Models (LLMs) represent a breakthrough in artificial intelligence, enabling machines to understand, generate, and process human-like text for tasks like answering questions, code generation, summarizing, and translating. LLM built on transformer architecture [13]. Mixture-of-Recursions (MoR) is a new model using parameter sharing, adaptive computation, and a routing mechanism to enhance the efficiency of language model.

This chapter presents reviews of existing studies related to Large Language Models (LLMs), Small Language Models (SLMs), transformer architectures, Recursive Transformer, the Mixture of Expert (MoE) and the Mixture-of-Recursions (MoR) framework. The goal is to identify the strengths and weaknesses of current models so that we will increase the effectiveness of our model.

### B. Large Language Model (LLM)

Large Language Models (LLMs) are language models trained on massive text corpora, predominantly utilizing the Transformer architecture. They are designed for natural language processing tasks, particularly text generation. By scaling parameters to billions and leveraging self-supervised learning, they acquire versatile capabilities in understanding and generating human-like text across diverse domains. The development of LLMs began with GPT-1, the first model to apply generative pre-training to a decoder-only Transformer [14].

While GPT-1 introduced generative pre-training, another model, GPT-3, a 175-billion parameter model, demonstrated that sufficiently large models could learn to follow instructions and solve diverse tasks with minimal examples, requiring no special training [3]. However, the empirical scaling approach of these early models lacked a theoretical foundation for optimal resource allocation, which was systematically addressed by Kaplan et al. [15]. They formulated scaling laws, demonstrating power-law relationships between model size, dataset size, training compute, and performance. Following these scaling laws, later models achieved even greater performance. Google’s PaLM used 540 billion parameters, effectively scaled across varied computational infrastructure, demonstrating breakthrough performance on reasoning tasks and achieving human-level proficiency [16]. GPT-4 outperforms human experts on numerous benchmarks, processing both text and image inputs to achieve human-level performance [17]. Despite these significant advancements, these models apply full model depth to every token, wasting computational resources on simple words.

### C. Small Language Model (SLM)

Small Language Models (SLMs) are Transformer-based generative models similar to Large Language Models (LLMs) but with significantly reduced dimensions. SLMs are lightweight versions of traditional language models; while LLMs typically possess 100 billion to over 10 trillion parameters, SLMs range from 1 million to 10 billion parameters. Although significantly smaller, they retain core NLP capabilities [18]. LLMs offer maximal general capability at high cost, whereas SLMs emphasize efficiency, privacy, and specialization.

However, the massive size, high cost, and complexity of LLMs pose significant challenges for certain applications, making them impractical for many use cases. SLMs were developed as a solution to these problems [19], aiming to make machine intelligence more accessible, affordable, and efficient for everyday tasks [20]. The SLM journey began in

2019 with DistilBERT, which compressed the BERT model from 110M to 66M parameters, achieving 97% of BERT’s performance with 60% faster inference, proving that small models could retain significant capability [21]. The modern SLM era expanded in 2023 with the Microsoft Phi family. Phi-1, at 1.3 billion parameters, achieved coding performance competitive with models over  $100\times$  larger [22]. Phi-2 achieved 58% on MMLU (Massive Multitask Language Understanding), outperforming 25B parameter LLMs [23]. Phi-3 Mini, with 3.8 billion parameters, reached 68.8% on MMLU, beating Llama 70B while running on smartphones [24]. This evolution proves that capability does not strictly require trillions of parameters.

However, a persistent inefficiency remains: all standard Transformer-based SLMs use the same number of layers for every token. This fixed-depth limitation creates an opportunity for MoR, which utilizes dynamic recursive depths to optimize computation.

#### D. Transformer

The Transformer architecture proposed by Vaswani et al. [12] introduces a mechanism known as “self-attention,” which replaces recurrence for sequence modeling, enabling parallel processing. This mechanism has dramatically improved training efficiency and handling of long-range dependencies compared to RNNs/LSTMs. Self-attention relates different positions of a single sequence to compute a representation of the sequence [12].

Transformers consist of layered encoder-decoder blocks containing multi-head attention and feed-forward sublayers. The decoder includes an additional attention layer that attends to the encoder stack outputs. Positional encoding allows the entire sequence to be processed in parallel while preserving word order information [12].

LLMs have revolutionized language modeling, with the Transformer architecture serving as their core foundation (e.g., BERT, GPT, T5). Despite successes in tasks like code generation and complex reasoning, this architecture faces limitations such as high computational cost, memory bottlenecks, and lack of adaptability. Bae et al. [25] introduce Relaxed Recursive Transformers utilizing parameter sharing to reduce model size. Building on this, the MoR (Mixture-of-Recursion) framework [1] combines parameter sharing with adaptive computation, integrating two key efficiency paradigms within a single architecture. This combined approach offers a potential solution to the fixed-depth limitations of standard Transformers.

#### E. Recursive Transformer

“Recursive Transformer” refers to architectures adding explicit recursion or iteration to standard Transformer blocks to improve efficiency or generalization. It emerged as a solution to the “fixed computational depth” limitation. The Universal Transformer introduced recurrence and adaptive computation time [26]. However, Universal Transformers do not function perfectly in practice, often performing well only on short sequences and applying uniform computation per word. Mixture-

of-Recursions (MoR) advances this by adding a “router” that evaluates each word individually [1].

#### F. Mixture of Experts

Mixture of Experts (MoE) is a neural network architecture designed to improve model capacity and efficiency using specialized subnetworks (“experts”) and a gating network that dynamically selects experts for each input. It involves scaling modern LLMs while keeping per-token compute constant [27], [28]. The concept originates from Jacobs et al. [29].

Shazeer et al. [11] introduced the Sparsely-Gated MoE layer for conditional computation, using top-k gating. To address training instability, Switch Transformers simplified this to routing each token to a single expert [30]. Jiang et al. [31] used two experts per token to improve balance. DeepSeek-V2 introduced further routing diversity [32]. MoR mixes routing strategies to optimize efficiency, specialization, and generalization.

#### G. Mixture of Recursion

Mixture-of-Recursion (MoR), introduced by Bae et al. [1], presents a unified architecture merging parameter sharing with adaptive computation. In contrast to standard Transformers where every input token is processed by a fixed sequence of layers [12], MoR employs a single shared layer stack that can be re-entered multiple times. Crucially, it assigns a specific processing depth to each token based on its complexity, allowing for variable compute paths.

**Parameter Sharing:** This technique reuses the weights of a layer across multiple recursive steps. It enables the construction of effectively deep networks without a corresponding increase in the number of distinct parameters, thereby reducing memory overhead [1], [33].

**Adaptive Computation:** MoR utilizes lightweight, learnable routers to dynamically determine the recursion depth for each token. This ensures that computational resources are concentrated on more difficult tokens that require deeper processing [26]. This approach differs from halting mechanisms, offering stable training dynamics by avoiding discrete, hard-to-optimize decisions [34].

**Recursion-Wise KV Cache and Routing:** To further optimize efficiency, MoR implements a specialized caching strategy where Key-Value (KV) pairs are stored only for tokens active at a specific recursion depth. The routing core employs a Multi-Layer Perceptron (MLP) to output probabilities for three distinct actions: Recurse (add depth), Forward (proceed), or Skip (bypass), effectively balancing performance and cost.

### III. METHODOLOGY

This chapter details the methodology employed to implement, train, and evaluate the MoR Transformer framework. The objective is to empirically assess the MoR architecture, originally proposed by Bae et al. [1], in the context of Small Language Models (SLMs). We investigate whether MoR’s dynamic depth adaptation confers efficiency benefits when training small-scale models from scratch. We benchmark MoR

Transformer efficiency against standard Transformers on a character-level language modeling task.

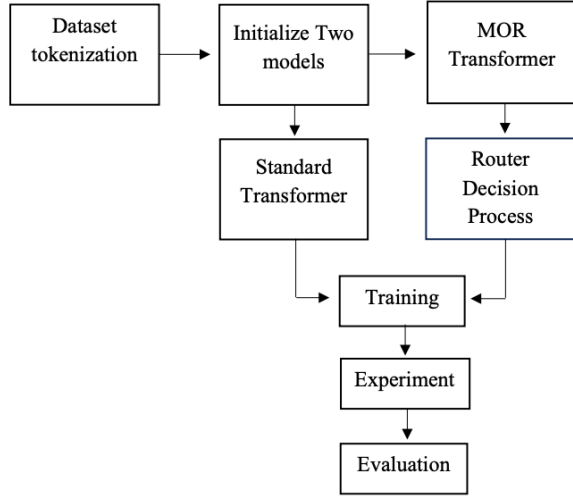


Fig. 1. Evaluated Research Pipeline

#### A. Dataset and Preprocessing

1) *Dataset*: The full Tiny Shakespeare dataset ( $\sim 1,000,000$  characters) is used as a controlled benchmark. It is sourced from the Andrej Karpathy charRNN project [35] and chosen for its linguistic structure and suitability for rapid experimentation. The vocabulary size is 55 characters.

2) *Tokenization*: Character-level tokenization is applied. Each character is mapped to a unique integer, as shown in Table I. The corpus is encoded as a one-dimensional tensor of token IDs.

TABLE I  
TOKENIZATION

Name	Description
stoi	Maps each character to a unique integer ID
itos	Inverse mapping from ID back to character

#### B. Research Pipeline

The research pipeline (Figure 1) operates sequentially: dataset tokenization, random initialization of MoR and Standard Transformer parameters, router integration, training from scratch, and finally, evaluation.

#### C. Model Architecture

1) *Baseline Transformer*: The baseline is a standard Transformer encoder-decoder language model [12]. It includes token embeddings, positional encoding, and a stack of  $N$  fixed layers.

2) *MoR-Transformer*: The MoR-Transformer uses the same backbone but introduces a routing mechanism. Each layer optionally includes a standard block and an MoR Router. The router decides to Skip, Execute, or Recurse.

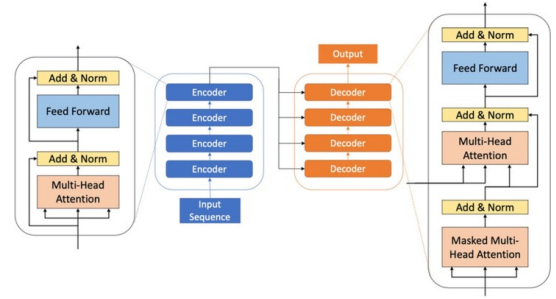


Fig. 2. Transformer Architecture Diagram

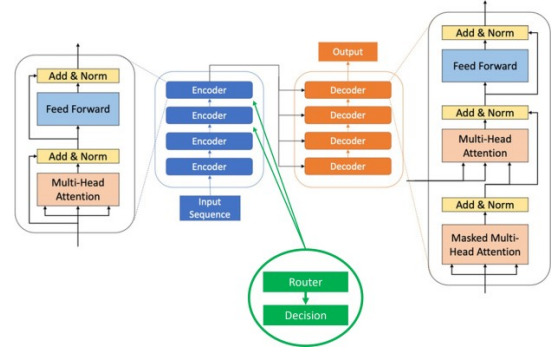


Fig. 3. MoR-Transformer Diagram

3) *Router Architecture*: The MoR router is a lightweight feed-forward network that determines the execution path for each token:

- **Input**: Token embeddings of dimension  $d_{model} = 256$
- **Architecture**: Two-layer MLP:  $256 \rightarrow 128 \rightarrow 3$  with ReLU activation
- **Output**: Softmax probability distribution over three actions: Skip (cost 0), Forward (cost 1), Recurse (cost 1+)
- **Training**: Gumbel-Softmax reparameterization with temperature  $\tau = 1.0$  for differentiable discrete sampling
- **Inference**: Argmax selection of highest-probability action
- **Regularization**: Auxiliary loss  $\lambda \cdot E$  penalizes effective depth to encourage efficient routing

#### D. Experiments

We conducted two experiments:

- **Experiment 1 (Efficiency Profiling)**: Compares a 12-layer Baseline with a 12-layer MoR-Transformer.
- **Experiment 2 (Equal Cost Comparison)**: Compares a 6-layer Baseline against a 12-layer MoR model (Effective depth  $E \approx 6$ ).

**Data Split**: Experiment 1 uses the full set for training/evaluation. Experiment 2 uses the first half for training and the second half as held-out test data.

#### E. Fine-tuning Process

Both models are trained from scratch on the full Tiny Shakespeare dataset. Configuration details are in Table II.

TABLE II  
FINE-TUNING CONFIGURATION

Parameter	Value/Description
Batch size	64
Architecture	Transformer/Mor-Transformer
d_model(hidden size)	256
Num. layers ( $N_{full}$ )	12
Num. layers baseline low ( $N_{low}$ )	6 (for equal-compute experiment)
Optimizer	AdamW (base params), AdamW (router params)
learning rate	1e-3
Baseline (all exp)	30 epochs
MoR (Exp 1)	30 epochs
MoR (Exp 2)	50 epochs (to optimize routing)
Platform and Hardware	Kaggle, Gpu

#### F. Evaluation metrics

To understand how well our models are working, we used some standard evaluation metrics. These help us measure accuracy and speed of our model.

1) *Effective Depth (E)*: Effective Depth (E) is the average number of layers actually executed by a model during inference or training. It is a measure of the realized computational depth of a dynamically-routed or recursive neural architecture.

2) *Accuracy*: Accuracy measures top-1 next-token prediction correctness for language modeling. It shows what fraction of predicted tokens match true next tokens.

3) *Held-out accuracy*: Held-out accuracy is the classification accuracy measured on data that was not used for training (the held-out or test set). It tells how well the model generalizes to unseen examples.

4) *Fine-tuning Time*: Fine-tuning time is the total amount of time a model spends to finish fine-tuning on a given dataset. It helps evaluate: Efficiency (how fast a model trains), Comparison between models (e.g., Standard Transformer vs. MoR), Compute cost and Time-to-convergence.

### IV. RESULT AND ANALYSIS

This chapter presents the experimental results obtained after training and evaluating the MoR and baseline models. The results are analyzed based on metrics such as accuracy, effective depth, and training time. The performance of the MoR-Transformer is discussed in detail to demonstrate its effectiveness.

#### A. Results

This section details the results of the Efficiency Profiling and Equal Cost Comparison experiments.

1) *Experiment 1: Efficiency Profiling ( $N=12$  vs MoR  $N=12$ )*: Goal: Compare a normal Transformer (12 layers executed always) with a MoR Transformer (12 recursive layers allowed to skip/recurse dynamically).

**Interpretation:** The results indicate that the MoR model achieves virtually identical predictive accuracy to the standard Transformer (approx. 22.76). Notably, the MoR router utilized an average of only 8 layers per sample, compared to the fixed 12 layers of the baseline. This 33% reduction in computations per token sequence underscores the efficiency gains of adaptive

TABLE III  
EXPERIMENT 1 RESULT

Metric	Standard (N=12)	MoR (N=12)
Accuracy	22.7634	22.7635
Effective Depth (E)	12.00	8.00
Fine-tuning Time	108 s	83 s

routing, leading to faster and more resource-efficient training without compromising accuracy.

2) *Experiment 2: Performance Under Equivalent Cost ( $N=6$  vs MoR  $N=12$ ,  $E \approx 6$ )*: Goal: See if a MoR model (12 potential layers, dynamic depth  $\approx 6$ ) can match the compute cost of a smaller baseline (6 layers) while achieving better accuracy through adaptive computation.

TABLE IV  
EXPERIMENT 2 RESULT

Metric	Standard (N=6)	MoR (N=12, $E \approx 6$ )
Held-out Accuracy	39.87 %	49.67 %
Effective Depth (E)	6.00	5.89
Fine-tuning Time	30 s	59 s

**Interpretation:** With an effective depth  $E \approx 6$ , the MoR model operated within the same average compute budget as the 6-layer baseline. Despite equivalent computational cost, the MoR model achieved a 10% absolute improvement in held-out accuracy, demonstrating superior generalization. The router effectively learned to deepen recursion for complex inputs while skipping layers for simpler ones, enabling the model to function as a variable-depth expert system that efficiently balances computation and accuracy.

TABLE V  
OVERALL OUTCOME

Experiment	What it proved	Key Takeaway
1: Efficiency	For same Layer (12 L), MoR reduces compute ( $E = 8$ ) with identical accuracy.	MoR learns to skip unnecessary computation.
2: Equivalent Cost	For same average cost ( $E \approx 6$ ), MoR (12 L) achieves higher accuracy than fixed 6 L.	MoR uses depth where needed $\rightarrow$ better adaptability and generalization.

#### 3) Overall Experimental Outcome:

#### B. Result Analysis

In Experiment 1, the MoR-Transformer reduced effective depth by 33% ( $E=8.00$  vs baseline  $E=12.00$ ) while keeping accuracy almost the same (22.7635 vs 22.7634). This shows that the routing mechanism can learn to skip unnecessary layers and control recursion under a strong auxiliary penalty ( $\lambda = 0.1$ ), which speeds up training by 23% (83s vs 108s) without performance degradation.

Experiment 2 revealed even stronger benefits: on held-out data, the MoR model ( $N=12$ ,  $E=5.89$ ) outperformed the shallower baseline ( $N=6$ ,  $E=6.00$ ) by 24.6% in accuracy (49.6687 vs 39.8700) at marginally lower computational cost,

confirming that MoR leverages its greater nominal capacity through adaptive execution paths—allocating deeper recursion to complex sequences while skipping simpler ones—thus validating both core hypotheses that MoR provides superior efficiency at equal depth and better generalization under equivalent compute budgets, consistent with broader MoR literature reporting 2x inference gains through dynamic recursion.

## V. CONCLUSION & FUTURE WORK

### A. Conclusion

This paper explored the effectiveness of the Mixture of Recursions (MoR) architecture as an adaptive computation mechanism for next-token prediction tasks, with a particular focus on small language models (SLMs). Unlike standard Transformers, which apply a fixed depth of computation to every input, MoR introduces a lightweight routing mechanism that allows each token to either skip, execute once, or recurse through a shared layer. This design enables the model to allocate more computation to complex tokens while saving resources on simpler ones.

Two experiments were conducted using a character-level dataset inspired by Tiny Shakespeare. The Experimental comparisons between MoR and standard Transformers demonstrate that MoR can achieve superior performance while using fewer effective layers. This model not only lowers effective depth but also preserves and improves accuracy compared to standard Transformers under the same computational budget. The study confirms MoR supports recursive reasoning, which is absent in traditional MoE models, and adapts its routing behavior based on token complexity. These capabilities make MoR a scalable and efficient solution for language modeling, especially in resource-constrained environments.

In conclusion, our findings position MoR as a viable alternative to fixed-depth Transformers, particularly in resource-constrained scenarios favoring SLMs, such as mobile and edge computing. By enhancing computational efficiency without sacrificing accuracy, MoR represents a significant step towards sustainable and scalable language modeling.

### B. Limitations

This study has several limitations that should be acknowledged:

- **Dataset Scale:** The Tiny Shakespeare dataset (~1M characters) is relatively small compared to modern language modeling benchmarks. Results may not generalize to larger, more diverse corpora.
- **Single Task:** Only character-level next-token prediction was evaluated. Performance on other tasks (e.g., code generation, question answering) remains unexplored.
- **Training from Scratch:** Models were trained from random initialization rather than fine-tuned from pre-trained weights, limiting insights into transfer learning scenarios.
- **Limited Baselines:** No comparison to other efficiency methods such as parameter-efficient fine-tuning (LoRA, Adapters), pruning, or early-exit mechanisms.

- **Statistical Rigor:** Results are based on single runs without multiple random seeds or significance testing.
- **Deployment Metrics:** Real-world inference latency, throughput, and energy consumption were not measured.

### C. Future Works

Although this research provides strong evidence for the benefits of adaptive recursion, several avenues remain open for future investigation. One important extension is to evaluate MoR on larger and more diverse datasets, including code generation, question answering, translation and summarization to better understand its generalization ability. Combine MoR with MoE to leverage both expert specialization and recursive depth control for enhanced performance. Finally, deploying MoR on resource constrained hardware like smartphones or IoT devices would demonstrate its practical benefits for realtime systems, and exploring its extension to multimodal tasks (text plus image or text plus code) could open new opportunities for adaptive recursion in broader domains.

## REFERENCES

- [1] S. Bae et al. Mixture-of-recursions: Learning dynamic recursive depths for adaptive token-level computation. In *Proc. 39th Int. Conf. Mach. Learn. (ICML)*, 2025.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] T. Brown, B. Mann, N. Ryder, et al. Language models are few-shot learners. In *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, volume 33, pages 1877–1901, 2020.
- [4] C. Raffel, N. Shazeer, A. Roberts, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1):5485–5551, 2020.
- [5] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proc. 57th Annu. Meet. Assoc. Comput. Linguist. (ACL)*, pages 2978–2988, 2019.
- [6] N. Kitaev, L. Kaiser, and A. Levskaya. Reformer: The efficient transformer. In *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [7] J. W. Rae, A. Potapenko, S. M. Jayakumar, C. Hillier, and T. P. Lillicrap. Compressive transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*, 2019.
- [8] L. Ben Allal et al. SmolLM2: When smol goes big—data-centric training of a small language model. *arXiv preprint*, 2025.
- [9] Y. Li, Y. Huang, M. E. Ildiz, A. S. Rawat, and S. Oymak. Mechanics of next token prediction with self-attention. *arXiv preprint*, 2024.
- [10] M. E. Sander and G. Peyré. Towards understanding the universality of transformers for next-token prediction. *arXiv preprint arXiv:2410.03011*, 2024.
- [11] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2017.
- [12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, pages 5998–6008, 2017.
- [13] H. Naveed et al. A comprehensive overview of large language models. *arXiv preprint*, 2023.
- [14] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. Improving language understanding by generative pre-training. Technical report, OpenAI, 2018.
- [15] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [16] A. Chowdhery et al. PaLM: Scaling language modeling with pathways. *J. Mach. Learn. Res.*, 24(240):1–113, 2023.
- [17] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

- [18] J. Okah. Small language models (SLM): A comprehensive overview. Hugging Face, 2025.
- [19] Q. Zhang et al. The rise of small language models. *IEEE Intell. Syst.*, 40(1):30–37, 2025.
- [20] Z. Lu et al. Small language models: Survey, measurements, and insights. *arXiv preprint*, 2025.
- [21] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [22] S. Gunasekar, Y. Zhang, J. Aneja, C. C. Teodoro Mendes, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- [23] M. Javaheripi et al. Phi-2: The surprising power of small language models. Technical report, Microsoft, 2023.
- [24] M. Abdin et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [25] S. Bae, A. Fisch, H. Harutyunyan, Z. Ji, S. Kim, and T. Schuster. Relaxed recursive transformers: Effective parameter sharing with layer-wise LoRA. In *Proc. 13th Int. Conf. Learn. Represent. (ICLR)*, 2025.
- [26] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and L. Kaiser. Universal transformers. In *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2019.
- [27] W. Cai et al. A survey on mixture of experts in large language models. *IEEE Trans. Knowl. Data Eng.*, 36(12):7304–7317, 2024.
- [28] Y. Zhou et al. Mixture-of-experts with expert choice routing. In *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, volume 35, pages 7103–7114, 2022.
- [29] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts. *Neural Comput.*, 3(1):79–87, 1991.
- [30] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23(1):5232–5270, 2022.
- [31] A. Q. Jiang et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [32] DeepSeek-AI. DeepSeek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024.
- [33] S. Shen et al. Sliced recursive transformer. *arXiv preprint*, 2022.
- [34] M. Elbayad, J. Gu, E. Grave, and M. Auli. Depth-adaptive transformer. In *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [35] A. Karpathy. char-rnn. GitHub repository, 2015. <https://github.com/karpathy/char-rnn>.